# Preconditioning Constrained Eigenvalue Problems[☆]

C.G. Baker[a], R.B. Lehoucq[b]

[a]*Scalable Algorithms Department, Sandia National Laboratories, Albuquerque, New Mexico 87185 USA*
[b]*Applied Mathematics and Applications Department, MS 1320, Sandia National Laboratories, Albuquerque, New Mexico 87185 USA*

## Abstract

The purpose of our paper is introduce a robust preconditioning scheme for the numerical solution of the constrained eigenvalue problem for approximating the leftmost eigenvalues and corresponding eigenvectors. This constrained eigenvalue problem is in general equivalent to a nonsymmetric eigenvalue problem with nontrivial Jordan blocks associated with infinite eigenvalues. The preconditioning scheme may be used in combination with Krylov subspace methods and preconditioned eigensolvers. The two key results are a semi-orthogonal decomposition and a transformation process that combines a preconditioning step followed by abstract projection onto the subspace associated with the finite eigenvalues. Numerical results demonstrate the effectiveness of the preconditioning scheme.

*Key words:* Eigenvalues, Eigenvectors, Preconditioning, Projection, linear equality constraints

## 1. Introduction

We are interested in computing the leftmost eigenvalues and corresponding eigenvectors of the constrained eigenvalue problem

$$\begin{bmatrix} \mathbf{A} & \mathbf{B}^T \\ \mathbf{B} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} = \begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} \lambda, \qquad \mathbf{A}, \mathbf{M} = \mathbf{M}^T \in \mathbb{R}^{n \times n}, \mathbf{B} \in \mathbb{R}^{m \times n}, m < n \tag{1}$$

where $\lambda \in \mathbb{C} \cup \{\infty\}$ is an eigenvalue, $\mathbf{u} \in \mathbb{C}^n$ and $\mathbf{v} \in \mathbb{C}^m$. We also assume that the saddle point matrix on the left hand side of (1) is invertible. In particular, we are interested "preconditioning" (1) so that the relative separation of the leftmost eigenvalues improve, and in exploiting a preconditioner $\mathbf{K}$ for $\mathbf{A}$. We also assume that only the application of $\mathbf{K}^{-1}$ upon a vector, and matrix vector products with $\mathbf{A}$ and $\mathbf{M}$ are available.

Computing approximations to the leftmost eigenvalues and corresponding eigenvectors of (1) is a nontrivial task when only iterative methods are employed. For instance, (1) contains $2m$ "infinite" eigenvalues so that critical to the success in computing the leftmost eigenvectors is that any approximation is a member of the subspace associated with finite eigenvalues. Preconditioning should not only improve the relative separation of the leftmost eigenvalues but avoid the eigenspace associated with the infinite eigenvalues. The purpose of our paper is to introduce a robust preconditioning scheme for the numerical solution of the constrained eigenvalue problem (1) for approximating the leftmost eigenvalues and corresponding eigenvectors.

Applications of interest include the incompressible Navier-Stokes equations [1], contact problems in linear elasticity [2], electromagnetics [3], and when the columns of $\mathbf{B}^T$ contain eigenvectors of the matrix pencil $(\mathbf{A}, \mathbf{M})$. The matrix $\mathbf{A}$ often represents the finite element discretization of a second order differential operator, $\mathbf{Bx} = \mathbf{0}$ is a matrix of constraints, and $\mathbf{M}$ is a mass matrix. The constrained eigenvalue problem (1) often arises as the optimality system for the constrained energy problem

$$\min_{\mathbf{x} \neq \mathbf{0}} \frac{1}{2} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{M} \mathbf{x}} \quad \text{subject to} \quad \mathbf{Bx} = \mathbf{0}, \qquad \mathbf{A} = \mathbf{A}^T, \mathbf{M} = \mathbf{M}^T \in \mathbb{R}^{n \times n}, \mathbf{B} \in \mathbb{R}^{m \times n}, m < n, \tag{2}$$

---

where $\mathbf{M}$ is positive definite, $\mathbf{A}$ is positive definite on the Null($\mathbf{B}$), and $\mathbf{B}$ is of full row rank. Without loss of generality, the primal and dual parts $\mathbf{u}$ and $\mathbf{v}$ of the eigenvector can be assumed to lie in $\mathbb{R}^n$ and $\mathbb{R}^m$, respectively. The dual vector corresponds to the vector of Lagrange multipliers used to enforce the constraints.

## 2. Structure of Eigenvalue problem

The generalized eigenvalue problem (1) is, in general, equivalent to standard nonsymmetric eigenvalue problem. This is true even if $\mathbf{A}$ is assumed to be a symmetric matrix. Because $\mathbf{M}$ is a symmetric positive definite matrix, there exists $\mu \in \mathbb{R}$ so that $\mathbf{A} - \mu\mathbf{M}$ is a symmetric positive definite matrix. Then

$$\begin{bmatrix} \mathbf{A} & \mathbf{B}^T \\ \mathbf{B} & \mathbf{0} \end{bmatrix} + \mu \begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{B}\,(\mathbf{A} - \mu\mathbf{M})^{-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A} - \mu\mathbf{M} & \mathbf{0} \\ \mathbf{0} & -\mathbf{B}\,(\mathbf{A} - \mu\mathbf{M})^{-1}\,\mathbf{B}^T \end{bmatrix} \begin{bmatrix} \mathbf{I} & (\mathbf{A} - \mu\mathbf{M})^{-1}\,\mathbf{B}^T \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$$

implies that the associated generalized symmetric eigenvalue problem (1) is indefinite because there are positive and negative eigenvalues. Hence (1) is equivalent to a nonsymmetric eigenvalue problem (see [4, Chapter 15]) regardless of whether $\mathbf{A}$ is a symmetric matrix. However, we can transform (1) into a symmetric positive semi-definite generalized eigenvalue problem. Subtract $\mu \begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$ from both sides of (1) to obtain

$$\begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{A} - \mu\mathbf{M} & \mathbf{B}^T \\ \mathbf{B} & \mathbf{0} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} = \begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} (\lambda - \mu)^{-1} \tag{3}$$

and note that the triple product of matrices above simplifies to

$$\begin{bmatrix} \mathbf{M}\,(\mathbf{A} - \mu\mathbf{M})^{-1}\,\mathbf{M} + \mathbf{H}^T\mathbf{H} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \qquad \mathbf{H} = (\mathbf{A} - \mu\mathbf{M})^{-1/2}\,\mathbf{B}^T\mathbf{B}\,(\mathbf{A} - \mu\mathbf{M})^{-1}\,\mathbf{M}$$

under the assumption that $\mathbf{A} - \mu\mathbf{M}$ is a symmetric positive definite matrix. Such a transformation allows the numerical solution of (3) by symmetric eigensolvers (assuming that any needed orthogonality is relaxed to semi-orthogonality with respect to the semi-inner product induced by $\begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$).

Denote the Rayleigh quotient of the vector $\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \in \mathbb{R}^{n+m}$ for the eigenvalue problem (1) by

$$\theta\left( \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \right) = \frac{\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}^T \begin{bmatrix} \mathbf{A} & \mathbf{B}^T \\ \mathbf{B} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}}{\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}^T \begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}} = \frac{\mathbf{x}^T\mathbf{A}\mathbf{x} + 2\mathbf{y}^T\mathbf{B}\mathbf{x}}{\mathbf{x}^T\mathbf{M}\mathbf{x}} \tag{4}$$

and the associated residual

$$\begin{bmatrix} \mathbf{A} & \mathbf{B}^T \\ \mathbf{B} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} - \begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \theta\left( \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \right) = \begin{bmatrix} \mathbf{r}_1 \\ \mathbf{r}_2 \end{bmatrix}. \tag{5}$$

Let $\mathbf{S} = \mathbf{S}^T \in \mathbb{R}^{n \times n}$ be a positive definite matrix, and let

$$\mathbf{P} = \mathbf{I} - \mathbf{S}^{-1}\mathbf{B}^T \left( \mathbf{B}\mathbf{S}^{-1}\mathbf{B}^T \right)^{-1} \mathbf{B}, \qquad \mathbf{Q} = \mathbf{I} - \mathbf{P}. \tag{6}$$

The matrices $\mathbf{P}$ and $\mathbf{Q}$ are $\mathbf{S}$ orthogonal projectors from $\mathbb{R}^n$ onto Null($\mathbf{B}$) and Range($\mathbf{B}^T$) in the direction of $\mathbf{S}^{-1}$Range($\mathbf{B}^T$), respectively. The following result provides an orthogonal decomposition that proves central to our considerations.

**Lemma 1.** *Let $\mathbf{B} \in \mathbb{R}^{m \times n}$ be of full row rank. For $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{y} \in \mathbb{R}^m$, the decomposition*

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{P}\mathbf{x} \\ \mathbf{w} \end{bmatrix} + \begin{bmatrix} \mathbf{0} \\ \mathbf{y} - \mathbf{w} \end{bmatrix} + \begin{bmatrix} \mathbf{Q}\mathbf{x} \\ \mathbf{0} \end{bmatrix} \tag{7}$$

*is orthogonal with respect to the semi-inner product induced by* $\begin{bmatrix} \mathbf{S} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{n+m \times n+m}$, *where* $\mathbf{w} \in \mathbb{R}^m$ *is given by*

$$\mathbf{B}\mathbf{S}^{-1}\mathbf{B}^T\mathbf{w} = \mathbf{B}\mathbf{S}^{-1}\left(\mathbf{M}\mathbf{P}\mathbf{x}\nu - \mathbf{A}\mathbf{P}\mathbf{x} + \mathbf{r}_1\right), \qquad (\mathbf{P}\mathbf{x})^T\,\mathbf{r}_1 = 0, \quad \mathbf{r}_2 = \mathbf{0} \tag{8}$$

*where* $\mathbf{r}_1$ *and* $\mathbf{r}_2$ *are given by* (5), *and*

$$\nu = \frac{(\mathbf{P}\mathbf{x})^T\,\mathbf{A}\mathbf{P}\mathbf{x}}{(\mathbf{P}\mathbf{x})^T\,\mathbf{M}\mathbf{P}\mathbf{x}}. \tag{9}$$

*Proof.* That the decomposition (7) is orthogonal with respect to the semi-inner product induced by $\begin{bmatrix} \mathbf{S} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$ follows from the properties of the $\mathbf{S}$ orthogonal projectors $\mathbf{P}$ and $\mathbf{Q}$. Because $\mathbf{B}\mathbf{P} = \mathbf{0}$, then (4) results in

$$\theta\left(\begin{bmatrix} \mathbf{P}\mathbf{x} \\ \mathbf{y} \end{bmatrix}\right) = \frac{(\mathbf{P}\mathbf{x})^T\,\mathbf{A}\mathbf{P}\mathbf{x}}{(\mathbf{P}\mathbf{x})^T\,\mathbf{M}\mathbf{P}\mathbf{x}} = \nu \in \mathbb{R}$$

and so (5) implies that $(\mathbf{P}\mathbf{x})^T\,\mathbf{r}_1 = 0$ and $\mathbf{r}_2 = \mathbf{0}$. If we premultiply both sides of (5) by $\begin{bmatrix} \mathbf{S}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$, then

$$\begin{cases} \mathbf{S}^{-1}\mathbf{B}^T\mathbf{w} &= \mathbf{S}^{-1}\left(\mathbf{M}\mathbf{P}\mathbf{x}\nu - \mathbf{A}\mathbf{P}\mathbf{x} + \mathbf{r}_1\right), \\ \mathbf{B}\mathbf{P}\mathbf{x} &= \mathbf{0}, \end{cases} \tag{10}$$

where we used the identity $\mathbf{B}\mathbf{P} = \mathbf{0}$ and $\mathbf{r}_1$ from (5). Equation (10) results in

$$\begin{cases} \mathbf{B}\mathbf{S}^{-1}\mathbf{B}^T\mathbf{w} &= \mathbf{B}\mathbf{S}^{-1}\left(\mathbf{M}\mathbf{P}\mathbf{x}\nu - \mathbf{A}\mathbf{P}\mathbf{x} + \mathbf{r}_1\right), \\ \mathbf{B}\mathbf{P}\mathbf{x} &= \mathbf{0}, \end{cases} \tag{11}$$

where $\mathbf{B}\mathbf{S}^{-1}\mathbf{B}^T \in \mathbb{R}^{m \times m}$ is a symmetric positive definite matrix. $\square$

We remark that if $\mathbf{S}$ is a nonsymmetric or symmetric indefinite matrix, then the decomposition is no longer orthogonal but the decomposition holds. Because $\mathbf{P}^2 = \mathbf{P}$, the matrix $\mathbf{P}$ is by definition a projector (albeit an *oblique* one). Denote

$$\mathcal{E} = \mathrm{Span}\left\{\begin{bmatrix} \mathbf{P}\mathbf{x} \\ \mathbf{w} \end{bmatrix}\right\}, \quad \mathcal{E}_\infty = \mathrm{Span}\left\{\begin{bmatrix} \mathbf{0} \\ \mathbf{c} \end{bmatrix}\right\}, \quad \mathcal{E}_d = \mathrm{Span}\left\{\begin{bmatrix} \mathbf{Q}\mathbf{x} \\ \mathbf{0} \end{bmatrix}\right\},$$

where $\mathbf{c} \in \mathbb{R}^m$. Now, premultiply both sides of (1) by $\begin{bmatrix} \mathbf{S}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$, and eliminate the dual variable $\mathbf{v}$ to obtain

$$\begin{cases} \mathbf{S}^{-1}\mathbf{P}^T\mathbf{A}\mathbf{u} &= \mathbf{S}^{-1}\mathbf{P}^T\mathbf{M}\mathbf{u}\lambda, \\ \mathbf{B}\mathbf{u} &= \mathbf{0}, \end{cases}$$

which can be rewritten as

$$\mathbf{S}^{-1}\mathbf{P}^T\mathbf{A}\mathbf{P}\mathbf{u} = \mathbf{S}^{-1}\mathbf{P}^T\mathbf{M}\mathbf{P}\mathbf{u}\lambda. \tag{12}$$

Premultiplying the previous equation by $(\mathbf{S}\mathbf{u})^T$ leads to

$$\lambda = \frac{(\mathbf{P}\mathbf{u})^T\,\mathbf{A}\mathbf{P}\mathbf{u}}{(\mathbf{P}\mathbf{u})^T\,\mathbf{M}\mathbf{P}\mathbf{u}} \in \mathbb{C}, \tag{13}$$

because of the assumption on $\mathbf{M}$ given by (1). The number of these finite eigenvalues is $n - m$ because the rank of $\mathbf{P}$ is $n - m$. Any vector in $\mathbf{z} \in \mathcal{E}_\infty$ leads to

$$\begin{bmatrix} \mathbf{A} & \mathbf{B}^T \\ \mathbf{B} & \mathbf{0} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{z} = \mathbf{z} \cdot 0,$$

where the matrix inverse exists due to the assumption on the saddle point system following (1). Hence, the dimension of $\mathcal{E}_\infty$ is $m$. The rank of $\mathbf{Q}$ is $m$ so that the dimension of $\mathcal{E}_d$ is $m$. The above discussion leads to the following result.

**Lemma 2.** *Let the hypothesis of Lemma 1 be given, and denote*

$$\mathcal{E} = \mathrm{Span}\left\{\begin{bmatrix}\mathbf{Px}\\\mathbf{w}\end{bmatrix}\right\}, \quad \mathcal{E}_\infty = \mathrm{Span}\left\{\begin{bmatrix}\mathbf{0}\\\mathbf{c}\end{bmatrix}\right\}, \quad \mathcal{E}_d = \mathrm{Span}\left\{\begin{bmatrix}\mathbf{Qx}\\\mathbf{0}\end{bmatrix}\right\}. \tag{14}$$

*Then*

$$\mathbb{R}^{n+m} = \mathcal{E} \oplus \mathcal{E}_\infty \oplus \mathcal{E}_d,$$

*is an* $\begin{bmatrix}\mathbf{S} & \mathbf{0}\\\mathbf{0} & \mathbf{0}\end{bmatrix}$ *semi-orthogonal decomposition, and*

1. $\mathcal{E}$ *is an eigenspace for* (1) *of dimension* $n - m$ *with complex eigenvalues,*
2. $\mathcal{E}_\infty$ *is an eigenspace for* (1) *of dimension* $m$ *associated with an infinite eigenvalue,*
3. $\mathcal{E}_d$ *is of dimension* $m$ *where* $\mathcal{E}_\infty$ *is orthogonal (with respect to the Euclidean inner product on* $\mathbb{R}^{n+m}$*) to* $\mathcal{E}_d$.

The above analysis emphasizes that the numerical solution of (1) depends upon the approximation to the finite eigenvectors remaining in $\mathcal{E}$. The usefulness of Lemmas 1 and 2 is that remaining in $\mathcal{E}$ is equivalent to maintaining orthogonality to the undesired space $\mathcal{E}_\infty \oplus \mathcal{E}_d$. Numerically, the computation is sensitive to the error in maintaining the desired orthogonality. As discussed after Lemma 1, the decomposition holds when $\mathbf{S}$ is a nonsymmetric or symmetric indefinite matrix but the decomposition can no longer be orthogonal. Any resulting numerical implementation is then unstable because the decomposition is bi-orthogonal.

Lemmas 1 and 2 generalize the decomposition introduced by Malkus [5] to an orthogonal one (let alone to an $\mathbf{S}$ orthogonal decomposition). Instead, Malkus provides a decomposition based on the Jordan Canonical form (see also Theorem 1 in [6]). Cliffe, Garratt, and Spence [1] present a decomposition $\mathbf{S} = \mathbf{I}$ using the QR factorization.

We end this section with the following result that proves useful for preconditioning (1). The Lemma describes a transformation mapping $\mathbb{R}^{n+m}$ to $\mathcal{E}$.

**Lemma 3.** *Let* $\mathbf{S} \in \mathbb{R}^{n\times n}$ *be a symmetric positive definite matrix, let* $\mathbf{C} \in \mathbb{R}^{n\times n}$, *and let* $\mathbf{x} \in \mathbb{R}^n$. *If* $\mathbf{P}$ *is given by* (6), *then*

$$\begin{bmatrix}\mathbf{S} & \mathbf{B}^T\\\mathbf{B} & \mathbf{0}\end{bmatrix}^{-1}\begin{bmatrix}\mathbf{C} & \mathbf{0}\\\mathbf{0} & \mathbf{0}\end{bmatrix}\begin{bmatrix}\mathbf{x}\\\mathbf{y}\end{bmatrix} = \begin{bmatrix}\mathbf{PS}^{-1}\mathbf{Cx}\\\left(\mathbf{BS}^{-1}\mathbf{B}^T\right)^{-1}\mathbf{BS}^{-1}\mathbf{Cx}\end{bmatrix} \in \mathcal{E} \tag{15}$$

*where* $\mathbf{y} \in \mathbb{R}^m$ *is arbitrary. In particular, if* $\mathbf{S}^{-1}\mathbf{Cx} \in \mathrm{Range}(\mathbf{P})$, *then*

$$\begin{bmatrix}\mathbf{S} & \mathbf{B}^T\\\mathbf{B} & \mathbf{0}\end{bmatrix}^{-1}\begin{bmatrix}\mathbf{C} & \mathbf{0}\\\mathbf{0} & \mathbf{0}\end{bmatrix}\begin{bmatrix}\mathbf{x}\\\mathbf{y}\end{bmatrix} = \begin{bmatrix}\mathbf{PS}^{-1}\mathbf{Cx}\\\mathbf{0}\end{bmatrix}. \tag{16}$$

*Proof.* A tedious application of the identity

$$\begin{bmatrix}\mathbf{S} & \mathbf{B}^T\\\mathbf{B} & \mathbf{0}\end{bmatrix}^{-1} = \begin{bmatrix}\mathbf{I} & -\mathbf{S}^{-1}\mathbf{B}^T\\\mathbf{0} & \mathbf{I}\end{bmatrix}\begin{bmatrix}\mathbf{S}^{-1} & \mathbf{0}\\\mathbf{0} & -\left(\mathbf{BS}^{-1}\mathbf{B}^T\right)^{-1}\end{bmatrix}\begin{bmatrix}\mathbf{I} & \mathbf{0}\\-\mathbf{BS}^{-1} & \mathbf{I}\end{bmatrix}$$

establishes the equality of (15). A simple calculation demonstrates that the vector on the righthand side of (15) satisfies (8) of Lemma 1 and so the membership of (15) follows. If $\mathbf{S}^{-1}\mathbf{Cx} \in \mathrm{Range}(\mathbf{P})$, then (16) follows because $\mathbf{BP} = \mathbf{0}$. $\qquad\square$

Note that the hypothesis on $\mathbf{S}$ can be weakened to a nonsymmetric invertible $\mathbf{S}$. The projector $\mathbf{P}$ is no longer orthogonal and so $\mathcal{E}$ is no longer semi-orthogonal to $\mathcal{E}_\infty \oplus \mathcal{E}_d$.

Lemma 3 also allows us to link the vector in $\mathcal{E}$ of (15) with the solution of the constrained minimization problem

$$\arg\min_{\mathbf{z}\in\mathbb{R}^n}\left(\frac{1}{2}\mathbf{z}^T\mathbf{Sz} - \mathbf{z}^T\mathbf{S}\left(\mathbf{S}^{-1}\mathbf{Cx}\right)\right) \quad \text{subject to} \quad \mathbf{Bz} = \mathbf{0}. \tag{17}$$

The optimality system for this constrained minimization problem is: Find $\hat{\mathbf{x}} \in \mathbb{R}^n$ and $\hat{\mathbf{y}} \in \mathbb{R}^m$ so that

$$\begin{bmatrix}\mathbf{S} & \mathbf{B}^T\\\mathbf{B} & \mathbf{0}\end{bmatrix}\begin{bmatrix}\hat{\mathbf{x}}\\\hat{\mathbf{y}}\end{bmatrix} = \begin{bmatrix}\mathbf{S}\left(\mathbf{S}^{-1}\mathbf{Cx}\right)\\\mathbf{0}\end{bmatrix} = \begin{bmatrix}\mathbf{C} & \mathbf{0}\\\mathbf{0} & \mathbf{0}\end{bmatrix}\begin{bmatrix}\mathbf{x}\\\mathbf{y}\end{bmatrix},$$

which has the same solution as (15). The vector $\hat{\mathbf{y}}$ contains the Lagrange multipliers needed to enforce the constraint $\mathbf{Bz} = \mathbf{0}$. In words, the vector $\hat{\mathbf{x}}$ is the $\mathbf{S}$ orthogonal projection of $\mathbf{S}^{-1}\mathbf{Cx}$ onto $\mathrm{Null}(\mathbf{B})$ in the direction of $\mathbf{S}^{-1}\mathrm{Range}(\mathbf{B}^T)$.

*2.1. Important Special Case*

Suppose that (1) is the optimality system for the constrained energy minimization (2) (so that **A** is a symmetric matrix positive definite on Null(**B**)). We may rewrite (13) as

$$\lambda = \frac{(\mathbf{Pu})^T \mathbf{APu}}{(\mathbf{Pu})^T \mathbf{MPu}} > 0, \tag{18}$$

where without loss of generality $\mathbf{u} \in \mathbb{R}^n$ and $\mathbf{v} \in \mathbb{R}^m$. Hence the $n - m$ finite eigenvalues associated with $\mathcal{E}$ in Lemma 3 are positive so that the leftmost eigenvalues are also the smallest ones. Moreover,

$$\min_{\mathbf{Pu} \neq \mathbf{0}} \frac{(\mathbf{Pu})^T \mathbf{APu}}{(\mathbf{Pu})^T \mathbf{MPu}} \leq \nu = \frac{(\mathbf{Px})^T \mathbf{APx}}{(\mathbf{Px})^T \mathbf{MPx}} \leq \max_{\mathbf{Pu} \neq \mathbf{0}} \frac{(\mathbf{Pu})^T \mathbf{APu}}{(\mathbf{Pu})^T \mathbf{MPu}}$$

because **Px** is a linear combination of the $n - m$ eigenvectors **Pu** of (12) associated with the positive eigenvalues.

It remains to relate the leftmost eigenvector of (1) and the solution to the constrained energy problem (2). This latter solution is given by

$$\min_{\mathbf{Px} \neq \mathbf{0}} \frac{(\mathbf{Px})^T \mathbf{APx}}{(\mathbf{Px})^T \mathbf{MPx}} = \min_{\mathbf{z} \in \mathcal{E}, \neq \mathbf{0}} \theta(\mathbf{z}) \tag{19}$$

where $\theta(\mathbf{z})$ is given by (4). In words, the eigenvector associated with the smallest eigenvalue of (1) solves the constrained energy problem (2). We also remark that by (19)

$$\min_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^T \mathbf{Ax}}{\mathbf{x}^T \mathbf{Mx}} \leq \theta(\mathbf{z}) \leq \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^T \mathbf{Ax}}{\mathbf{x}^T \mathbf{Mx}} \qquad \mathbf{z} \in \mathcal{E} \oplus \mathcal{E}_d, \mathbf{z} \neq \mathbf{0}$$

so that an error in the direction of $\mathcal{E}_d$ leads to a violation of the constraints. [1] The Rayleigh quotient $\theta(\mathbf{z})$ may then be smaller than the minimizing energy associated with (2). In contrast, from (4) and $|\mathbf{y}^T \mathbf{Bx}| \leq \|\mathbf{B}\| \, \|\mathbf{x}\| \, \|\mathbf{y}\|$ we have that when $\mathbf{z} \in \mathbb{R}^{n+m}$

$$\min_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^T \mathbf{Ax}}{\mathbf{x}^T \mathbf{Mx}} - 2\|\mathbf{B}\| \frac{\|\mathbf{x}\|}{\|\mathbf{M}^{1/2}\mathbf{x}\|} \frac{\|\mathbf{y}\|}{\|\mathbf{M}^{1/2}\mathbf{x}\|} \leq \theta(\mathbf{z}) \leq \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^T \mathbf{Ax}}{\mathbf{x}^T \mathbf{Mx}} + 2\|\mathbf{B}\| \frac{\|\mathbf{x}\|}{\|\mathbf{M}^{1/2}\mathbf{x}\|} \frac{\|\mathbf{y}\|}{\|\mathbf{M}^{1/2}\mathbf{x}\|}.$$

Then $\theta(\mathbf{z})$ is unbounded as $\|\mathbf{y}\|$ increases relative to $\|\mathbf{M}^{1/2}\mathbf{x}\|$ so that an error in the direction of $\mathcal{E}_\infty$ can lead to an extremely small or large Rayleigh quotient. The culprit is that the denominator of Rayleigh quotient $\theta(\cdot)$ is unaffected by **y**. Therefore errors in the direction of $\mathcal{E}_\infty \oplus \mathcal{E}_d$ lead to unbounded under- and over-estimates of the smallest eigenvalue of (1).

# 3. Preconditioning the constrained eigenvalue problem

This section explains how Lemma 3 effects preconditioning by 1) considering a shift-invert spectral transformation and 2) how to exploit a preconditioner **K** for **A**. Both 1) and 2) precondition (1) by improving the relative separation of the leftmost eigenvalues. An important practical consideration is that Lemma 3 explains that we can:

**Option 1** solve a saddle point linear system with coefficient matrix $\begin{bmatrix} \mathbf{S} & \mathbf{B}^T \\ \mathbf{B} & \mathbf{0} \end{bmatrix}$, or

**Option 2** solve the linear set of equation $\mathbf{S}\hat{\mathbf{x}} = \mathbf{Cx}$ followed by the projector vector product $\mathbf{P}\hat{\mathbf{x}}$.

The saddle point system can be solved with a sparse direct solver, or with a preconditioned iterative method (see [7, 8, 9] for overviews and citations to the literature). Success of a preconditioned iterative method largely depends upon the quality of the preconditioner—often a nontrivial task that depends upon **S** and **B**. Conversely, the efficient implementation of **Option 2** depends upon the ability to solve efficiently linear set of equations with **S** followed by application of the projector **P** on the vector $\hat{\mathbf{x}}$. This latter step is equivalent to **S** orthogonalizing $\hat{\mathbf{x}}$ against the columns of $\mathbf{B}^T$.

---

[1] These bounds are sharp. For instance, select the columns of $\mathbf{B}^T$ to contain the $n - 1$ eigenvectors orthogonal to the largest eigenvector of the matrix pencil $(\mathbf{A}, \mathbf{M})$.

## 3.1. Shift-invert transformation

A conventional approach to solve (1) is to use a shift-invert spectral transformation. The generalized eigenvalue problem (1) is recast as the standard eigenvalue problem

$$\begin{bmatrix} \mathbf{A} - \sigma\mathbf{M} & \mathbf{B}^T \\ \mathbf{B} & \mathbf{0} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} = \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} (\lambda - \sigma)^{-1} \qquad \sigma \in \mathbb{R}. \tag{20}$$

The infinite eigenvalues of (1) are thus transformed to zero eigenvalues under the shift-invert spectral transformation. Hence the shift-invert spectral transformation *preconditions* (1) by improving the relative separation of the leftmost eigenvalues. If $\sigma$ is chosen near the smallest eigenvalue of (1), then this eigenvalue is large in magnitude under the spectral transformation. If $\mathbf{A} - \sigma\mathbf{M}$ is positive definite, then by Lemma 3, the vector $\begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix}(\lambda - \sigma)^{-1} \in \mathcal{E}$. Hence the spectral transformation has the interesting interpretation as effecting the $\mathbf{A} - \sigma\mathbf{M}$ orthogonal projection of $(\mathbf{A} - \sigma\mathbf{M})^{-1}\mathbf{Mx}$ onto Null($\mathbf{B}$) in the direction of $(\mathbf{A} - \sigma\mathbf{M})^{-1}$ Range($\mathbf{B}^T$). Lemma 3 implies that the Krylov space

$$\mathcal{K}_m\left(\begin{bmatrix} \mathbf{A} - \sigma\mathbf{M} & \mathbf{B}^T \\ \mathbf{B} & \mathbf{0} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \mathbf{z}\right) \in \mathcal{E}, \qquad m \geq 0 \tag{21}$$

if $\mathbf{z} \in \mathcal{E}$. This is accomplished by application of the shift-invert matrix (given by matrix product on the left hand side of (20)) to a vector in $\mathbb{R}^{n+m}$. Because the shift-invert matrix is nonsymmetric, the Arnoldi algorithm is employed. By the discussion following Lemma 3, nonsymmetric $\mathbf{A}$ leads to an oblique projector $\mathbf{P}$.

Premultiplication of (20) by $\begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$ results in the generalized eigenvalue problem (3) where $\mu = \sigma$. In contrast to (21), however, Lemma 3 implies that the Krylov space

$$\mathcal{K}_m\left(\begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}\begin{bmatrix} \mathbf{A} - \sigma\mathbf{M} & \mathbf{B}^T \\ \mathbf{B} & \mathbf{0} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \mathbf{z}\right) \in \mathbb{R}^{n+m}, \qquad m \geq 0 \tag{22}$$

even if $\mathbf{z} \in \mathcal{E}$. In practise, a basis is constructed for (22) by the Arnoldi algorithm by computing an $\begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$ semi-orthogonal basis for the Krylov subspace (21). If $\mathbf{A}$ is symmetric, then the shift-invert Lanczos method [10] can be applied because (3) defines a symmetric positive semi-definite generalized eigenvalue problem as discussed following (3). Both [6, 10] explain how a "purification" can be performed in an implicit fashion to remove the component in the direction of $\mathcal{E}_\infty \oplus \mathcal{E}_d$. By Lemma 3 this purification is equivalent to projection onto Null($\mathbf{B}$) via a shift-invert transformation.

We conclude this subsection that care must be taken when preconditioned iterative methods are used for **Option 1** and **Option 2** discussed at the start of §3. The tolerances for terminating the necessary solves and/or applications of the projectors must result in errors in the directions of the subspace $\mathcal{E}_d \oplus \mathcal{E}_\infty$ that are sufficiently small so as not to contaminate the computations. The norm of the application of $\mathbf{B}$ on approximations to the eigenvectors computed should be monitored.

## 3.2. Preconditioned Eigensolvers

We now introduce a preconditioning scheme for (1) that avoids the need to apply a shift-invert transformation as discussed in subsection 3.1. The results of [11] suggest that preconditioned eigensolvers can have a significant impact for eigenvalue problems arising in structural dynamics. Such preconditioned eigensolvers include gradient-based methods DACG [12, 13] and LOBPCG [14], the Davidson-based methods [15] such as the Newton based Jacobi-Davidson [16] algorithms, RTR [17], and trace minimization methods [18, 19]. Given an approximation $\mathbf{z} \in \mathcal{E}$ and the Rayleigh quotient $\theta(\mathbf{z})$, we apply a preconditioner $\mathbf{N} \in \mathbb{R}^{n+m \times n+m}$ to the residual $\mathbf{N}^{-1}\begin{bmatrix} \mathbf{r}_1 \\ \mathbf{r}_2 \end{bmatrix}$ where $\mathbf{r}_1$ and $\mathbf{r}_2$ are given by (5). In general, though, the preconditioned residual is a member of $\mathbb{R}^{n+m}$. However, once again, Lemma 3 explains how judicious choices for $\mathbf{S}$ and $\mathbf{C}$ map the preconditioned residual to $\mathcal{E}$. For instance, select

$$\mathbf{S} \leftarrow \mathbf{M}^{-1}, \qquad \mathbf{C} \leftarrow \mathbf{K}^{-1} \tag{23}$$

where $\mathbf{K}$ is a preconditioner for $\mathbf{A}$. Hence we see that judicious choices for $\mathbf{S}$ and $\mathbf{C}$ are equivalent *first* preconditioning the residual $\mathbf{r}_1$ followed by projection of $\begin{bmatrix} \mathbf{K}^{-1}\mathbf{r}_1 \\ \mathbf{0} \end{bmatrix}$ onto $\mathcal{E}$.

**Option 1** discussed at the start of §3 requires the solution of the saddle point linear system

$$\begin{bmatrix} \mathbf{M}^{-1} & \mathbf{B}^T \\ \mathbf{B} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{x}} \\ \hat{\mathbf{y}} \end{bmatrix} = \begin{bmatrix} \mathbf{K}^{-1}\mathbf{r}_1 \\ \mathbf{0} \end{bmatrix} \tag{24}$$

suggesting that $\mathbf{S} \leftarrow \mathbf{I}$ is a more expedient selection if the saddle point solver requires matrix-vector products with $\mathbf{M}^{-1}$. On the other hand **Option 2** leads to the following procedure

$$\begin{array}{lll} \text{Compute residual:} & \mathbf{r}_1 \\ \text{Precondition residual:} & \mathbf{K}^{-1}\mathbf{r}_1 \\ \text{Matrix-vector product:} & \hat{\mathbf{x}} \leftarrow \mathbf{M}\mathbf{K}^{-1}\mathbf{r}_1 \\ \text{Project:} & \mathbf{P}\hat{\mathbf{x}} \end{array} \quad . \tag{25}$$

Once again, as explained at the end of subsection 3.1 the tolerances for terminating the necessary solves and/or applications of the projectors must result in errors in the directions of the subspace $\mathcal{E}_d \oplus \mathcal{E}_\infty$ that are sufficiently small so as not to contaminate the computations. The norm of the application of $\mathbf{B}$ on approximations to the eigenvectors computed should be monitored.

## 4. Numerical experiments

We consider two eigenvalue problems. Problem 1 uses the symmetric positive definite, and constraint matrices

$$\mathbf{A}_1 = n \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & -1 & \ddots & & \\ & & & 2 & -1 \\ & & & -1 & 1 \end{bmatrix}, \quad \mathbf{M} = \frac{1}{6n} \begin{bmatrix} 4 & 1 & & & \\ 1 & 4 & 1 & & \\ & 1 & \ddots & & \\ & & & 4 & 1 \\ & & & 1 & 2 \end{bmatrix}, \quad \mathbf{B}^T = \mathbf{e}_{n/2} + \mathbf{e}_{n/2+1} \tag{26}$$

where $\mathbf{e}_i$ is the column $i$ of the identity matrix of order $n$ (assumed to be an even integer). Problem 2 replaces $\mathbf{A}_1$ with the nonsymmetric matrix

$$\mathbf{A}_2 = n \begin{bmatrix} 2 & -1-\alpha & & & \\ -1+\alpha & 2 & -1-\alpha & & \\ & -1+\alpha & \ddots & & \\ & & & 2 & -1-\alpha \\ & & & -1+\alpha & 1 \end{bmatrix} \tag{27}$$

where $\alpha$ is chosen small enough so that the matrix pencil $(\mathbf{A}_2, \mathbf{M})$ has real eigenvalues and eigenvectors.

We employ the restarted Davidson method described in [11] using MATLAB to compute the leftmost eigenvalue and eigenvector of (1). Orthogonality of the Davidson basis is maintained using the semi-inner product induced by the matrix $\begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$. All our experiments restart the Davidson method when the number of basis vectors constructed is 10, and are initialized with the same random vector in $\mathcal{E}$.

We first demonstrate that a naive application (so ignoring Lemma 3) of the Davidson method to solve 1 using the matrices (26), where $n = 100$, leads to a failure because of contamination by components in the direction $\mathcal{E}_\infty$ and $\mathcal{E}_d$. The residual calculated at each step of the Davidson iteration is preconditioned by the application of the block diagonal matrix $\begin{bmatrix} \mathbf{R}^T\mathbf{R} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{r}_1 \\ \mathbf{r}_2 \end{bmatrix}$ where $\mathbf{R}^T\mathbf{R}$ is the incomplete Cholesky of $\mathbf{A}_1$ formed via MATLAB's `cholinc(A,0.1)`, and $\mathbf{r}_1$ and $\mathbf{r}_2$ are given by (5). Figure 1 plots the Rayleigh quotient, residual norm, and constraint satisfaction of the Davidson iteration. The primal component of the initial iterate satisfies the constraint, a result of the initialization of
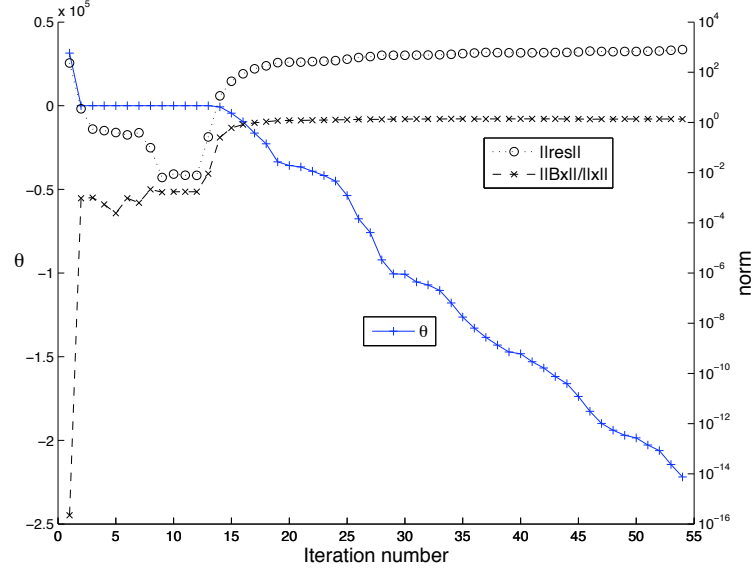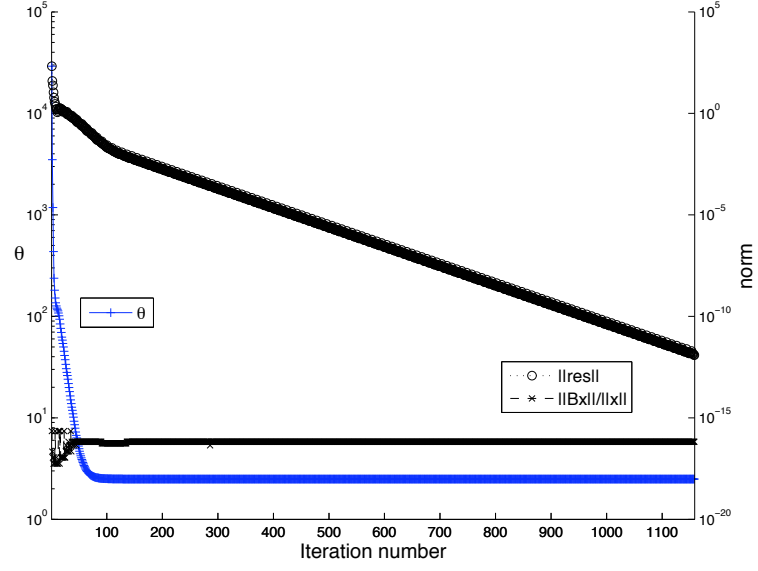
Figure 1: Davidson iteration for problem 26, where $n = 100$, with block diagonal preconditioning of the residual $[\mathbf{r}_1^T \ \mathbf{r}_2^T]^T$.
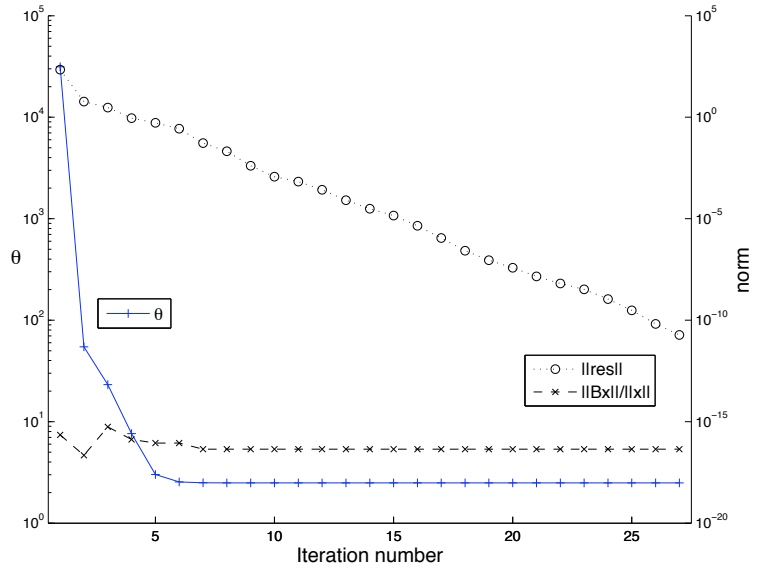
| Davidson iteration with preconditioner $\mathbf{K}_1$ | | | | | | |
|---|---|---|---|---|---|---|
| Tolerance $\epsilon$ | 1e-02 | 1e-04 | 1e-06 | 1e-08 | 1e-10 | 1e-12 |
| $\left\|\begin{bmatrix}\mathbf{r}_1\\\mathbf{r}_2\end{bmatrix}\right\|_2$ | 1.1e-02 | 1.2e-04 | 1.1e-06 | 1.1e-08 | 1.1e-10 | 1.5e-12 |
| $\|\mathbf{P}^T\mathbf{AP}\mathbf{x} - \mathbf{P}^T\mathbf{MP}\mathbf{x}\nu\|_2$ | 1.1e-02 | 1.2e-04 | 1.1e-06 | 1.1e-08 | 1.1e-10 | 1.5e-12 |
| $\theta$ | 2.49739 | 2.49192 | 2.49192 | 2.49192 | 2.49192 | 2.49192 |
| $\nu$ | 2.49739 | 2.49192 | 2.49192 | 2.49192 | 2.49192 | 2.49192 |
| $\|\mathbf{Bx}\|_2/\|\mathbf{x}\|_2$ | 5.5e-17 | 6.6e-17 | 6.6e-17 | 6.6e-17 | 6.6e-17 | 6.6e-17 |
| Davidson iteration with preconditioner $\mathbf{K}_2$ | | | | | | |
| Tolerance $\epsilon$ | 1e-02 | 1e-04 | 1e-06 | 1e-08 | 1e-10 | 1e-12 |
| $\left\|\begin{bmatrix}\mathbf{r}_1\\\mathbf{r}_2\end{bmatrix}\right\|_2$ | 5.9e-03 | 1.2e-04 | 4.0e-07 | 1.0e-08 | 9.5e-11 | 2.7e-11 |
| $\|\mathbf{P}^T\mathbf{AP}\mathbf{x} - \mathbf{P}^T\mathbf{MP}\mathbf{x}\nu\|_2$ | 5.9e-03 | 1.2e-04 | 4.0e-07 | 1.0e-08 | 9.5e-11 | 2.7e-11 |
| $\theta$ | 2.49193 | 2.49192 | 2.49192 | 2.49192 | 2.49192 | 2.49192 |
| $\nu$ | 2.49193 | 2.49192 | 2.49192 | 2.49192 | 2.49192 | 2.49192 |
| $\|\mathbf{Bx}\|_2/\|\mathbf{x}\|_2$ | 4.4e-17 | 4.4e-17 | 4.4e-17 | 4.4e-17 | 4.4e-17 | 4.4e-17 |

Table 1: Table of Rayleigh quotients, the norms of the residual and errors in the constraint are listed for two different choices of preconditioner for the problems associated with subfigures 2(a) and 2(b). The tolerance $\epsilon$ used to terminate the Davidson iterations is varied. $\theta$ is the Rayleigh quotient from (4) and $\nu$ is the Rayleigh quotient of ($\mathbf{P}^T\mathbf{AP}, \mathbf{P}^T\mathbf{MP}$) from (9), and $\|\mathbf{Bx}\|_2/\|\mathbf{x}\|_2$ measures the error in the constraint.

8

(a) Davidson iteration with preconditioner $\mathbf{K}_1$.



(b) Davidson iteration with preconditioner $\mathbf{K}_2$.

Figure 2: Davidson iteration for problem (26) where $n = 100$. The approximation to the smallest eigenvalue, the norms of the residual and errors in the constraint are displayed for two different choices of preconditioner.

the iteration with a vector in $\mathcal{E}$. However, as the Davidson iteration attempts convergence, the iterate becomes contaminated with components in $\mathcal{E}_\infty$ and $\mathcal{E}_d$. The Rayleigh quotient moves towards negative infinity, and the Davidson iteration fails.

Our next two examples exploit Lemma 3 within the Davidson iteration using the choices for $\mathbf{S}$ and $\mathbf{C}$ given by (23). The two examples use two different preconditioners $\mathbf{K}$ for $\mathbf{A}_1$ given by

$$\mathbf{K}_1 = \mathbf{I}, \qquad \mathbf{K}_2 = \mathbf{R}^T\mathbf{R}$$

where $\mathbf{R}^T\mathbf{R}$ is again the incomplete Cholesky factorization of $\mathbf{A}_1$ formed via MATLAB's `cholinc(A,0.1)`. We use **Option 1** as in (24) to apply Lemma 3. The Davidson method admits convergence when the normalized residual (5), using the semi-inner product induced by the matrix $\begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$, is less than the convergence tolerance $\epsilon$.

Figure 2 and Table 1 provide results for the two choices of preconditioners $\mathbf{K}_1$ and $\mathbf{K}_2$ and varying tolerances $\epsilon$. Because (12) can also be used to compute the finite eigenvalues and eigenvectors of (1), Table 1 also lists the norm of the residual, and the eigenvalue estimate $\nu$ (given by (9)). As a result of Lemma 3 and the starting conditions, both choices of preconditioner lead to a Davidson subspace in $\mathcal{E}$ for the duration of the iteration. Furthermore, the incomplete Cholesky factorization $\mathbf{K}_2 = \mathbf{R}^T\mathbf{R}$ of $\mathbf{A}_1$ improves the convergence rate of the Davidson method by reducing the number of iterations to 30 from 1100 (FIgure 2) when no preconditioning is used (or $\mathbf{K}_1 = \mathbf{I}$). We emphasize that the choice of preconditioner does not affect the eventual solution, only the rate of convergence.
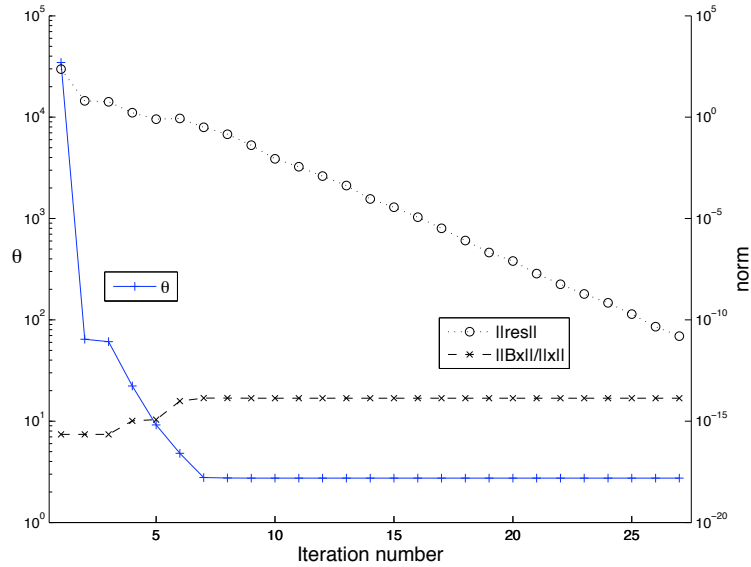


Figure 3: Davidson iteration for (27), where $n = 100$ using preconditioner $\mathbf{K}_3$.

In our last example, we apply the preconditioned Davidson method (modified for nonsymmetric eigenvalue problems) to the solution of the non-symmetric constrained eigenvalue problem given by (27). The preconditioner used is

$$\mathbf{K}_3 = \mathbf{LU},$$

an incomplete LU factorization of $\mathbf{A}_2$, and again we use the choices for $\mathbf{S}$ and $\mathbf{C}$ given by (23), and we use **Option 1** as in (24) to apply Lemma 3. These choices satisfy the conditions of Lemma 3 so that the Davidson subspace remains in $\mathcal{E}$. Figure 4 illustrates the successful trajectory of the Davidson iteration.

### References

[1] K. A. Cliffe, T. J. Garratt, A. Spence, Eigenvalues of block matrices arising from problems in fluid mechanics, SIAM J. Matrix Anal. Appl. 15 (1994) 1310–1318.

[2] R. D. Cook, D. S. Malkus, M. E. Plesha, R. J. Witt, Concepts and Applications of Finite Element Analysis, 3rd Edition, John Wiley & Sons, 2007.

[3] P. Arbenz, R. Geus, A comparison of solvers for large eigenvalue problems occuring in the design of resonant cavities, Numer. Linear Algebra Appl. 6 (1999) 3–16. `doi:10.1002/(SICI)1099-1506(199901/02)6:1<3::AID-NLA142>3.0.CO;2-I`.

[4] B. N. Parlett, The symmetric eigenvalue problem, SIAM, 1998.

[5] D. Malkus, Eigenproblems associated with the discrete LBB condition for incompressible finite elements, Int. J. Eng. Sci. 19 (1981) 1299–1310.

[6] K. Meerbergen, A. Spence, Implicitly restarted Arnoldi and purification for the shift-invert transformation, Mathematics of Computation 66 (1997) 667–689.

[7] M. Benzi, G. Golub, J. Liesen, Numerical solution of saddle point problems, Acta Numerica 14 (2005) 1–137.

[8] H. C. Elman, D. J. Silvester, A. J. Wathen, Finite Elements and Fast Iterative Solvers with applications in incompressible fluid dynamics, Numerical Mathematics and Scientific Computation, Oxford University Press, 2005.

[9] P. S. Vassilevski, Multilevel Block Factorization Preconditioners: Matrix-based Analysis and Algorithms for Solving Finite Element Equations, 1st Edition, Springer, 2008. `doi:10.1007/978-0-387-71564-3`.

[10] B. Nour-Omid, B. N. Parlett, T. Ericsson, P. Jensen, How to implement the spectral transformation, Mathematics of Computation 48 (1987) 663–673.

[11] P. Arbenz, U. L. Hetmaniuk, R. B. Lehoucq, R. S. Tuminaro, A comparison of eigensolvers for large-scale 3d modal analysis using AMG-preconditioned iterative methods, International Journal for Numerical Methods in Engineering 64 (2) (2005) 204–236. `doi:10.1002/nme.1365`.

[12] G. Gambolati, F. Sartoretto, P. Florian, An orthogonal accelerated deflation technique for large symmetric eigenvalue problem, Comput. Methods Appl. Mech. Engrg. 94 (1992) 13–23.

[13] L. Bergamaschi, G. Pini, F. Sartoretto, Approximate inverse preconditioning in the parallel solution of sparse eigenproblems, Numer. Linear Algebra Appl. 7 (3) (2000) 99–116.

[14] A. Knyazev, Toward the optimal preconditioned eigensolver: Locally optimal block preconditioned conjugate gradient method, SIAM J. Sci. Comput. 23 (2001) 517–541.

[15] E. R. Davidson, The iterative calculation of a few of the lowest eigenvalues and corresponding eigenvectors of large real-symmetric matrices, J. Comput. Phys. 17 (1975) 87–94.

[16] G. L. G. Sleijpen, H. A. van der Vorst, A Jacobi-Davidson iteration method for linear eigenvalue problems, SIAM J. Matrix Anal. Appl. 17 (2) (1996) 401–425.

[17] P.-A. Absil, C. G. Baker, K. A. Gallivan, Trust-region methods on Riemannian manifolds, Found. Comput. Math. 7 (3) (2007) 303–330. `doi:10.1007/s10208-005-0179-9`.

[18] A. Sameh, Z. Tong, The trace minimization method for the symmetric generalized eigenvalue problem, J. Comput. Appl. Math. 123 (2000) 155–175.

[19] A. H. Sameh, J. A. Wisniewski, A trace minimization algorithm for the generalized eigenvalue problem, SIAM J. Numer. Anal. 19 (6) (1982) 1243–1259.